



Text Processing using Machine Learning

Deep Learning Foundations + Nuts and Bolts

Liling Tan

23 Jan 2019

OVER
5,500 GRADUATE
ALUMNI

OFFERING OVER
120 ENTERPRISE IT, INNOVATION
& LEADERSHIP PROGRAMMES

TRAINING OVER
120,000 DIGITAL LEADERS
& PROFESSIONALS

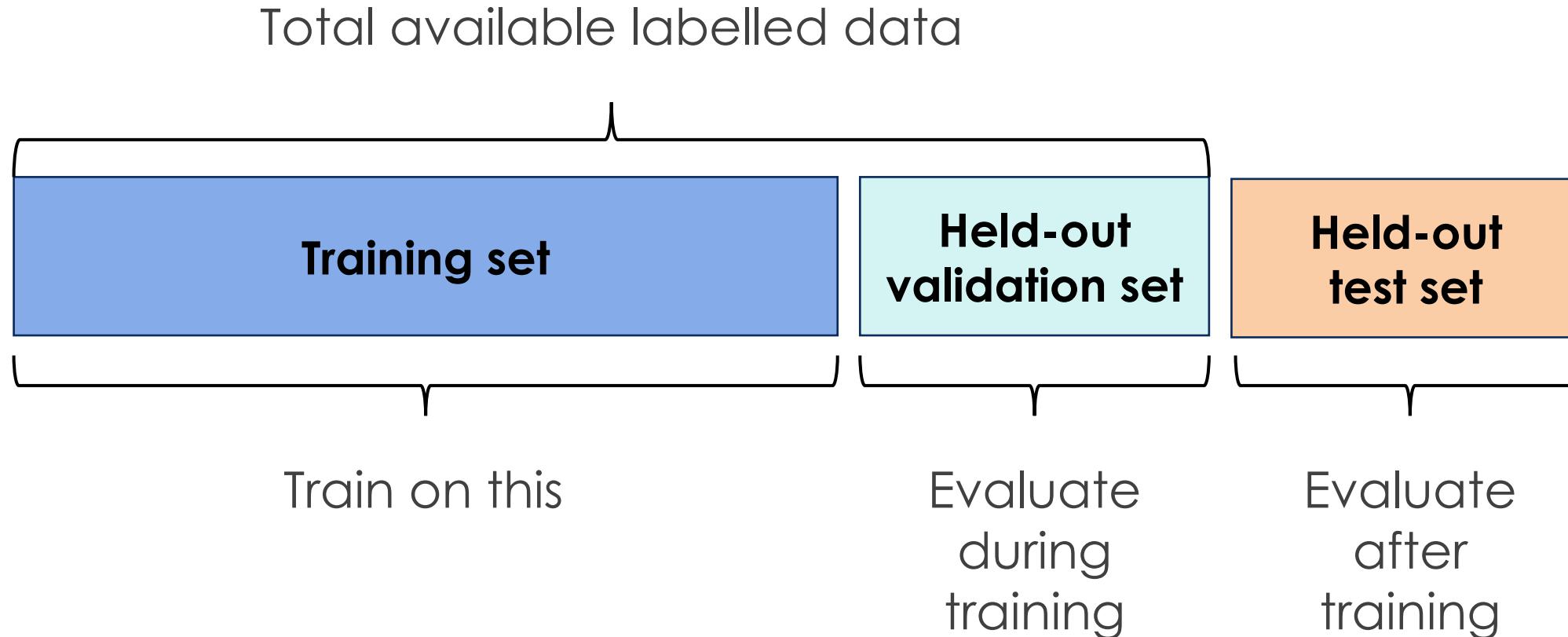
Overview

Lecture

- **Word2Vec Catchup**
 - Skipgram (20 mins)
 - Pretrained Embeddings (15 mins)
- **Nuts and Bolts**
 - Data Splits (5 mins)
 - Bias and Variance (10 mins)
 - Overfitting (10 mins)
 - Last Activation and Loss Function (15 mins)
 - Optimizers (15 mins)
- **TensorboardX** (60 mins)

Nuts and Bolts

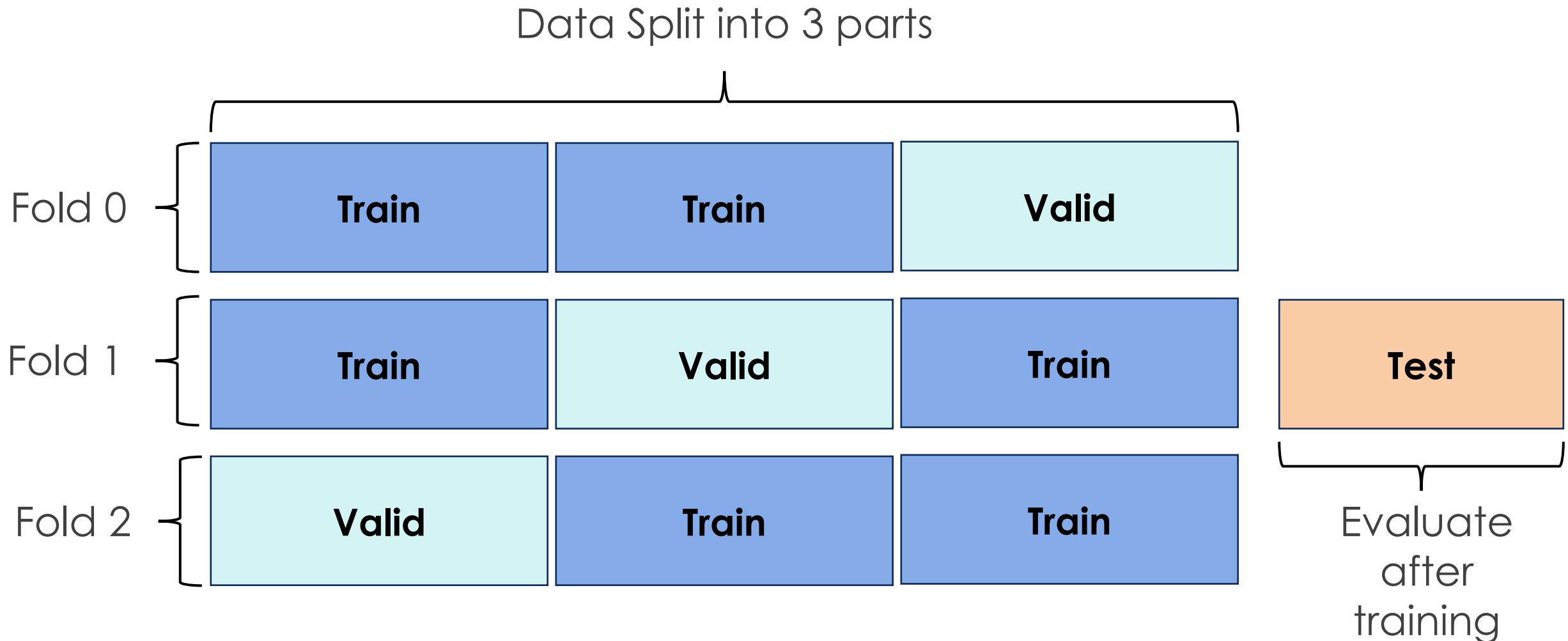
Hold-Out Evaluation



Hold-Out Evaluation

```
3 >>> from sklearn.model_selection import train_test_split
4 >>> import torch
5
6 >>> x, y = torch.rand(10, 2).numpy(), torch.rand(10).numpy()
7 >>> print(x.shape, y.shape)
8 (10, 2) (10,)
9
10 >>> split = 0.2
11 >>> x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=split)
12 >>> x_train, x_valid, y_train, y_valid = train_test_split(x_train, y_train, test_size=split)
13
14 >>> print(x_train.shape, x_valid.shape, x_test.shape)
15 (6, 2) (2, 2) (2, 2)
```

K-Fold Cross Validation Evaluation



K-Fold Cross Validation Evaluation

```
3  >>> from sklearn.model_selection import train_test_split
4  >>> from sklearn.model_selection import KFold
5  >>> import torch
6
7  >>> x, y = torch.rand(10, 2).numpy(), torch.rand(10).numpy()
8  >>> print(x.shape, y.shape)
9  (10, 2) (10,)
10
11 # Split out the held-out test set first.
12 >>> split = 0.2
13 >>> x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=split)
14
15 >>> kf = KFold(n_splits=3)
16 >>> three_folds = {i:{'x_train': x_train[train_idx], 'x_valid': x_train[valid_idx],
17 ...                           'y_train': y_train[train_idx], 'y_valid': y_train[valid_idx]}
18 ...           for i, (train_idx, valid_idx) in enumerate(kf.split(x_train))}
```

K-Fold Cross Validation Evaluation

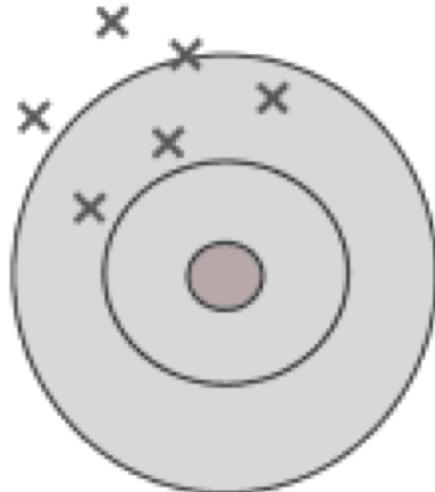
```
20 >>> three_folds[0]['x_train'].shape # Fold 0 train set.  
21 (5, 2)  
22 >>> three_folds[1]['x_train'].shape # Fold 1 train set.  
23 (5, 2)  
24 >>> three_folds[2]['x_train'].shape # Fold 2 train set.  
25 (6, 2)  
26  
27 >>> three_folds[0]['x_valid'].shape # Fold 0 valid set.  
28 (3, 2)  
29 >>> three_folds[1]['x_valid'].shape # Fold 1 valid set.  
30 (3, 2)  
31 >>> three_folds[2]['x_valid'].shape # Fold 2 valid set.  
32 (2, 2)
```

Bias-Variance Tradeoff

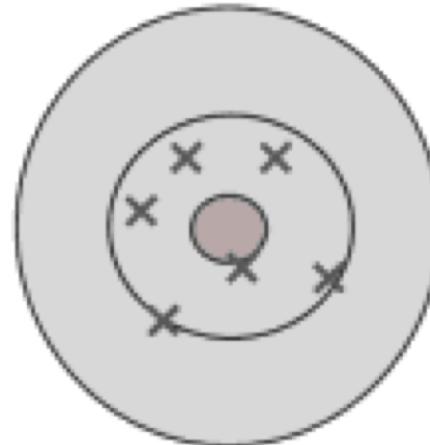
- “A **small network**, with say one hidden unit **is likely to be biased**, since the repertoire of available functions spanned by $f(x, w)$ over allowable weights will in this case be quite limited.”
- “if we **overparameterize**, via a large number of hidden units and associated weights, then bias will be reduced (... **with enough weights and hidden units, the network will interpolate the data**) but there is then the **danger of significant variance** contribution to the mean-square error”

(German et al. 1992)

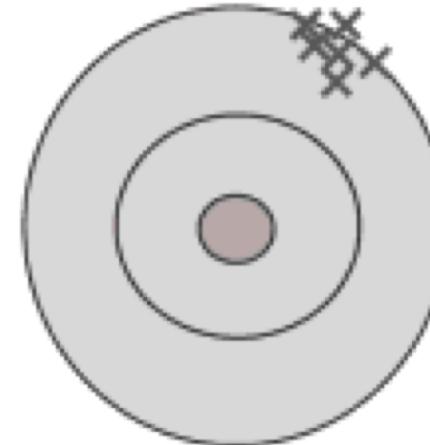
Bias-Variance Tradeoff



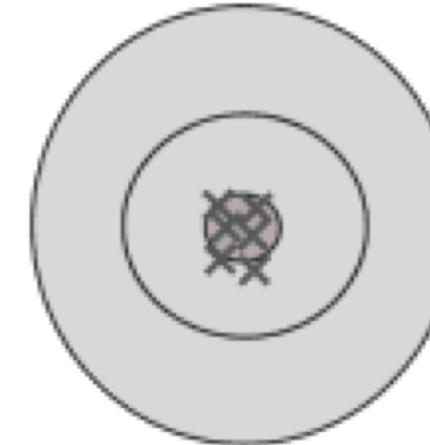
High bias
High variance



Low bias
High variance



High bias
Low variance



Low bias
Low variance

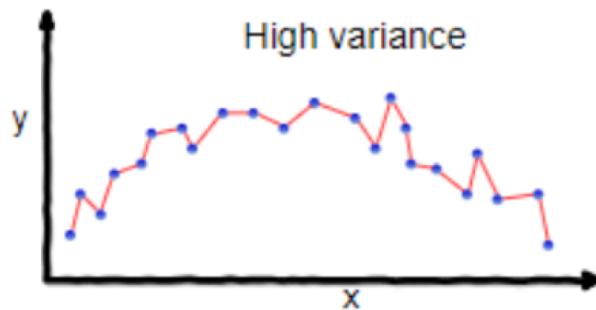
(Moore and McCabe, 2009)

Optimization and Generalization

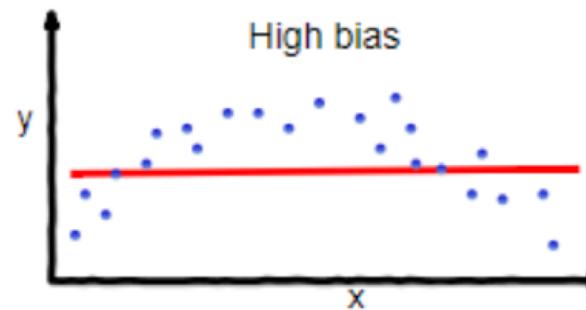
- **Optimization:** process of adjusting a model to get the best performance possible on the training data
- **Generalization:** how well trained model performs on data it has never seen before

Over-/Underfitting

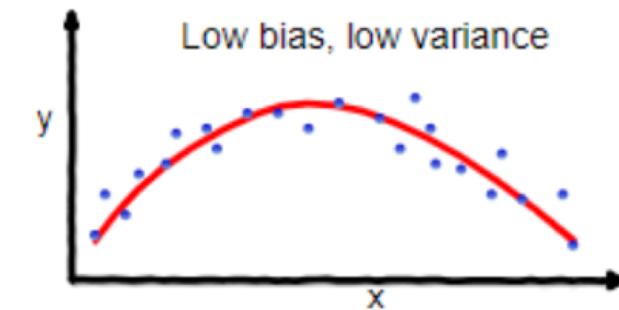
- **Underfit:** when model is not optimized
- **Overfit:** when model *fails to generalize*



overfitting



underfitting



Good balance

Optimization and Generalization

- **Optimization:** process of adjusting a model to get the best performance possible on the training data
- **Generalization:** how well trained model performs on data it has never seen before
- **What to do when model is underfitting (not optimal)?**
 - Train longer
- **How to prevent overfitting (i.e. generalize)?**
 - Get more data
 - Modulate quantity of information fed to model

Overcoming Overfitting

- **Reducing size of the model** (i.e. no. of layers and no. of units per layer) prevents overfitting
- **Weights regularization** (i.e. adding a cost associated with having large weights) put constraints on complexity of the model by forcing its weights to take small values
- **Dropping out** (i.e. randomly setting activated outputs to zero) is effective in regularizing the model



Last Layer Activation & Lost Functions

Creates a criterion that measures the mean squared error (squared L2 norm) between each element in the input x and target y .

The loss can be described as:

$$\ell(x, y) = L = \{l_1, \dots, l_N\}^\top, \quad l_n = (x_n - y_n)^2,$$

where N is the batch size. If `reduce` is `True`, then:

$$\ell(x, y) = \begin{cases} \text{mean}(L), & \text{if } \text{size_average} = \text{True}, \\ \text{sum}(L), & \text{if } \text{size_average} = \text{False}. \end{cases}$$

The sum operation still operates over all the elements, and divides by n .

The division by n can be avoided if one sets `size_average` to `False`.

To get a batch of losses, a loss per batch element, set `reduce` to `False`. These losses are not averaged and are not affected by `size_average`.

The negative log likelihood loss. It is useful to train a classification problem with C classes.

If provided, the optional argument *weight* should be a 1D Tensor assigning weight to each of the classes. This is particularly useful when you have an unbalanced training set.

The input given through a forward call is expected to contain log-probabilities of each class. *input* has to be a Tensor of size either (*minibatch*, C) or (*minibatch*, C , d_1, d_2, \dots, d_K) with $K \geq 2$ for the K -dimensional case (described later).

Obtaining log-probabilities in a neural network is easily achieved by adding a *LogSoftmax* layer in the last layer of your network. You may use *CrossEntropyLoss* instead, if you prefer not to add an extra layer.

The target that this loss expects is a class index (0 to $C-1$, where $C = \text{number of classes}$)

If `reduce` is `False`, the loss can be described as:

$$\ell(x, y) = L = \{l_1, \dots, l_N\}^\top, \quad l_n = -w_{y_n} x_{n,y_n}, \quad w_c = \text{weight}[c] \cdot 1\{c \neq \text{ignore_index}\},$$

where N is the batch size. If `reduce` is `True` (default), then

$$\ell(x, y) = \begin{cases} \sum_{n=1}^N \frac{1}{\sum_{n=1}^N w_{y_n}} l_n, & \text{if } \text{size_average} = \text{True}, \\ \sum_{n=1}^N l_n, & \text{if } \text{size_average} = \text{False}. \end{cases}$$

Can also be used for higher dimension inputs, such as 2D images, by providing an input of size $(minibatch, C, d_1, d_2, \dots, d_K)$ with $K \geq 2$, where K is the number of dimensions, and a target of appropriate shape (see below). In the case of images, it computes NLL loss per-pixel.

CrossEntropyLoss

This criterion combines `nn.LogSoftmax()` and `nn.NLLLoss()` in one single class.

It is useful when training a classification problem with C classes. If provided, the optional argument `weight` should be a 1D *Tensor* assigning weight to each of the classes. This is particularly useful when you have an unbalanced training set.

The *input* is expected to contain scores for each class.

input has to be a *Tensor* of size either $(minibatch, C)$ or $(minibatch, C, d_1, d_2, \dots, d_K)$ with $K \geq 2$ for the K -dimensional case (described later).

This criterion expects a class index (0 to $C-1$) as the *target* for each value of a 1D tensor of size *minibatch*

CrossEntropyLoss

The loss can be described as:

$$\text{loss}(x, \text{class}) = -\log \left(\frac{\exp(x[\text{class}])}{\sum_j \exp(x[j])} \right) = -x[\text{class}] + \log \left(\sum_j \exp(x[j]) \right)$$

or in the case of the *weight* argument being specified:

$$\text{loss}(x, \text{class}) = \text{weight}[\text{class}] \left(-x[\text{class}] + \log \left(\sum_j \exp(x[j]) \right) \right)$$

The losses are averaged across observations for each minibatch.

Can also be used for higher dimension inputs, such as 2D images, by providing an input of size $(\text{minibatch}, C, d_1, d_2, \dots, d_K)$ with $K \geq 2$, where K is the number of dimensions, and a target of appropriate shape (see below).

Creates a criterion that measures the Binary Cross Entropy between the target and the output:

The loss can be described as:

$$\ell(x, y) = L = \{l_1, \dots, l_N\}^\top, \quad l_n = -w_n [y_n \cdot \log x_n + (1 - y_n) \cdot \log(1 - x_n)],$$

where N is the batch size. If reduce is `True`, then

$$\ell(x, y) = \begin{cases} \text{mean}(L), & \text{if size_average} = \text{True}, \\ \text{sum}(L), & \text{if size_average} = \text{False}. \end{cases}$$

This is used for measuring the error of a reconstruction in for example an auto-encoder. Note that the targets y should be numbers between 0 and 1.

Last Layer Activation and Loss Function

Problem	Last Layer Activation	Loss Function	PyTorch
Binary Classification			
Multi-class, single-label classification			
Multi-class, multi-label classification			
Regression to arbitrary value			
Regression (0, 1)			

Last Layer Activation and Loss Function

Problem	Last Layer Activation	Loss Function	PyTorch
Binary Classification	Sigmoid	Binary Cross Entropy	<code>torch.nn.BCELoss</code>
Multi-class, single-label classification			
Multi-class, multi-label classification			
Regression to arbitrary value			
Regression (0, 1)			

Last Layer Activation and Loss Function

Problem	Last Layer Activation	Loss Function	PyTorch
Binary Classification	Sigmoid	Binary Cross Entropy	<code>torch.nn.BCELoss</code>
Multi-class, single-label classification	Sigmoid / Softmax		
Multi-class, multi-label classification			
Regression to arbitrary value			
Regression (0, 1)			

Last Layer Activation and Loss Function

Problem	Last Layer Activation	Loss Function	PyTorch
Binary Classification	Sigmoid	Binary Cross Entropy	<code>torch.nn.BCELoss</code>
Multi-class, single-label classification	Sigmoid / Softmax	Categorical Cross Entropy	<code>torch.nn.CrossEntropyLoss</code>
	LogSoftmax		
Multi-class, multi-label classification			
Regression to arbitrary value			
Regression (0, 1)			

Last Layer Activation and Loss Function

Problem	Last Layer Activation	Loss Function	PyTorch
Binary Classification	Sigmoid	Binary Cross Entropy	<code>torch.nn.BCELoss</code>
Multi-class, single-label classification	Sigmoid / Softmax	Categorical Cross Entropy	<code>torch.nn.CrossEntropyLoss</code>
	LogSoftmax	Negative Log Loss	<code>torch.nn.NLLLoss</code>
Multi-class, multi-label classification			
Regression to arbitrary value			
Regression (0, 1)			

Last Layer Activation and Loss Function

Problem	Last Layer Activation	Loss Function	PyTorch
Binary Classification	Sigmoid	Binary Cross Entropy	<code>torch.nn.BCELoss</code>
Multi-class, single-label classification	Sigmoid / Softmax	Categorical Cross Entropy	<code>torch.nn.CrossEntropyLoss</code>
	LogSoftmax	Negative Log Loss	<code>torch.nn.NLLLoss</code>
Multi-class, multi-label classification	Sigmoid / Softmax	Binary Cross Entropy	<code>torch.nn.BCELoss</code>
Regression to arbitrary value			
Regression (0, 1)			

Last Layer Activation and Loss Function

Problem	Last Layer Activation	Loss Function	PyTorch
Binary Classification	Sigmoid	Binary Cross Entropy	<code>torch.nn.BCELoss</code>
Multi-class, single-label classification	Sigmoid / Softmax	Categorical Cross Entropy	<code>torch.nn.CrossEntropyLoss</code>
	LogSoftmax	Negative Log Loss	<code>torch.nn.NLLLoss</code>
Multi-class, multi-label classification	Sigmoid / Softmax	Binary Cross Entropy	<code>torch.nn.BCELoss</code>
Regression to arbitrary value	None	L2 Loss	<code>torch.nn.MSELoss</code>
Regression (0, 1)	Sigmoid	L2 Loss	<code>torch.nn.MSELoss</code>

Optimizers

Gradient Descents

Stochastic Gradient Descent (SGD) $\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i)}; y^{(i)})$.
torch.optim.SGD

Mini-batch SGD
torch.optim.SGD

SGD with momentum
torch.optim.SGD

Gradient Descents

```
def step(self, closure=None):
    """Performs a single optimization step.

    Arguments:
        closure (callable, optional): A closure that reevaluates the model
            and returns the loss.
    """
    loss = None
    if closure is not None:
        loss = closure()

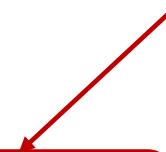
    for group in self.param_groups:
        weight_decay = group['weight_decay']
        momentum = group['momentum']
        dampening = group['dampening']
        nesterov = group['nesterov']
```

Gradient Descents

```
for p in group['params']:
    if p.grad is None:
        continue
    d_p = p.grad.data
    if weight_decay != 0:
        d_p.add_(weight_decay, p.data)
    if momentum != 0:
        param_state = self.state[p]
        if 'momentum_buffer' not in param_state:
            buf = param_state['momentum_buffer'] = torch.zeros_like(p.data)
            buf.mul_(momentum).add_(d_p)
        else:
            buf = param_state['momentum_buffer']
            buf.mul_(momentum).add_(1 - dampening, d_p)
        if nesterov:
            d_p = d_p.add(momentum, buf)
        else:
            d_p = buf
    p.data.add_(-group['lr'], d_p)

return loss
```

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$



Adaptive Moment Estimation ([Kingma, 2015](#))

- “momentum can be seen as a ball running down a slope, Adam behaves like a heavy ball with friction, which thus prefers flat minima in the error surface” – ([Dozat, 2016](#))

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t.$$

SGD vs Adam

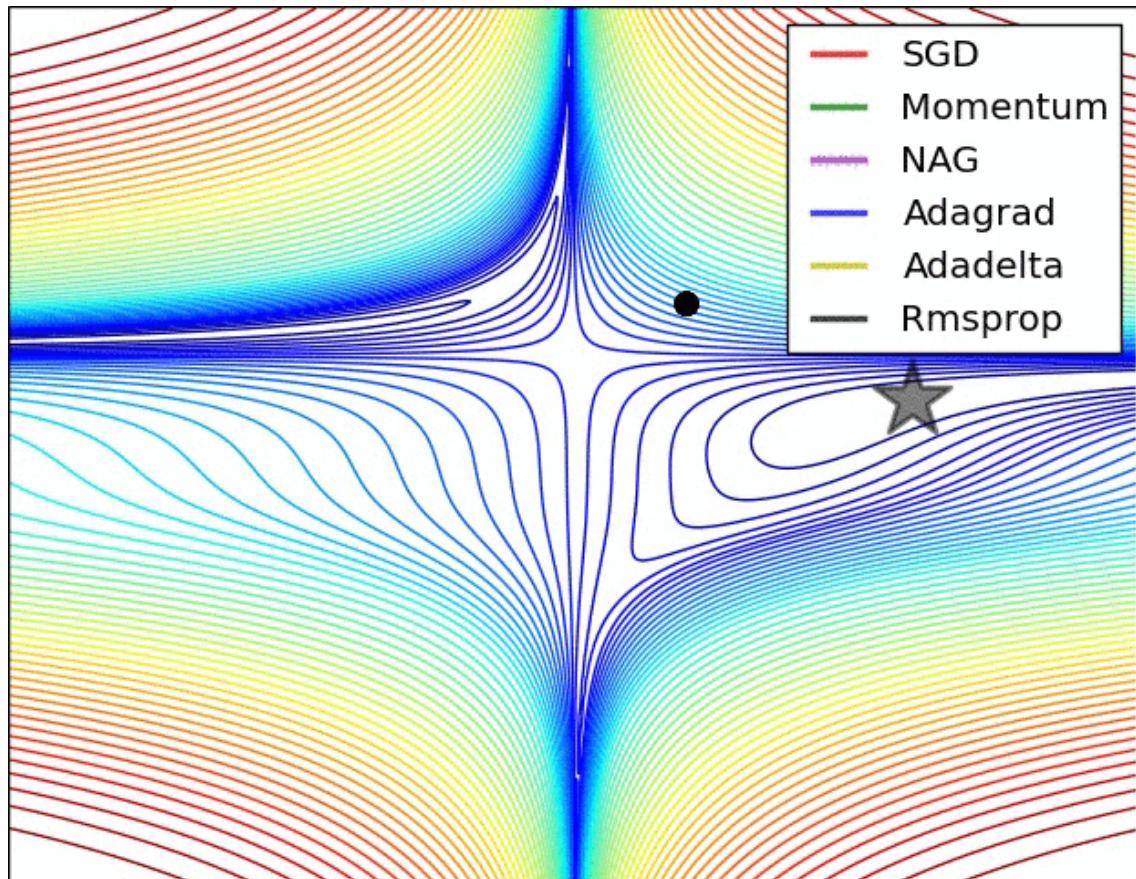
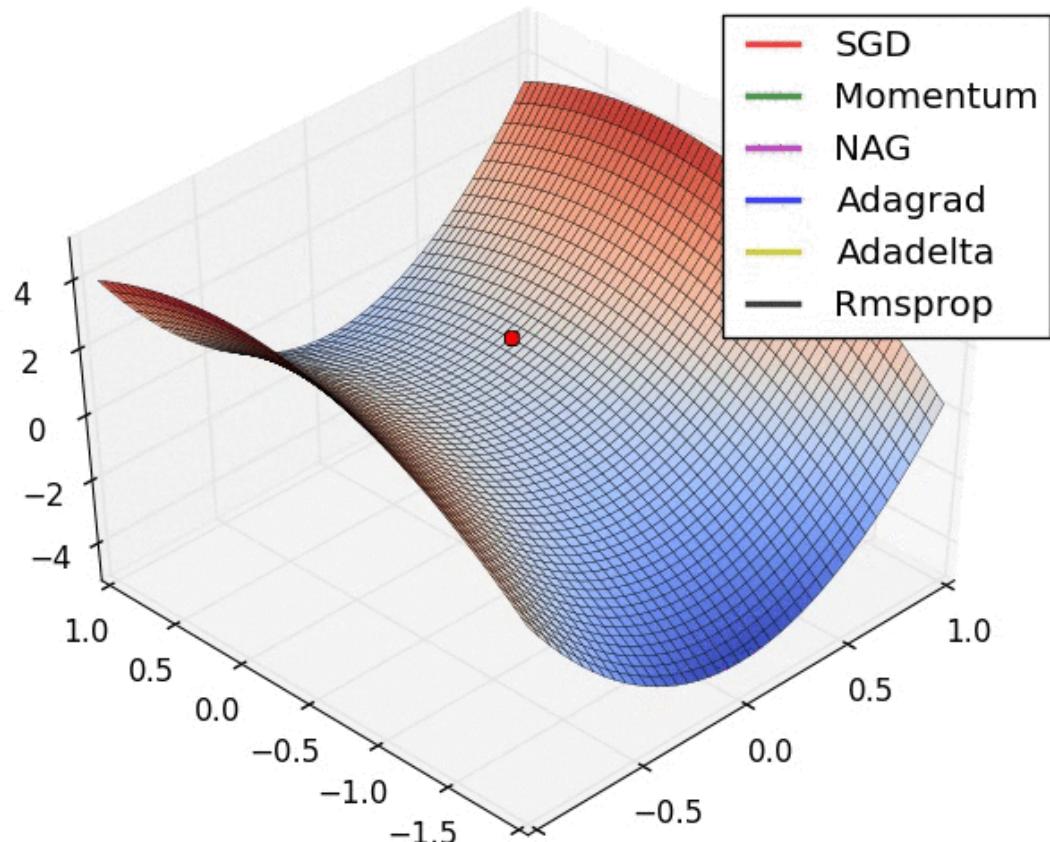
```
# Vanilla SGD
x += - learning_rate * dx
```

x is a vector of parameters and
dx is the gradient

```
# Adam
m = beta1*m + (1-beta1)*dx
v = beta2*v + (1-beta2)*(dx**2)
x += - learning_rate * m /
(np.sqrt(v) + eps)
```

x is a vector of parameters and
dx is the gradient
m is the smoothen gradient
v is the ‘cache’ used to normalize **x**
eps is smoothing term (1e-4 to 1e-8)
beta1, beta2 are hypers (0.9, 0.999)

Optimizers (Effects)



Note:

Deniz Yuret has a nice blogpost on optimizers <http://www.denizyuret.com/2015/03/alec-radfords-animations-for.html>

Sebastian Ruder too on optimizers <http://ruder.io/optimizing-gradient-descent/>

Stanford has a nice overview of optimizers on <http://cs231n.github.io/neural-networks-3/>

Google says...

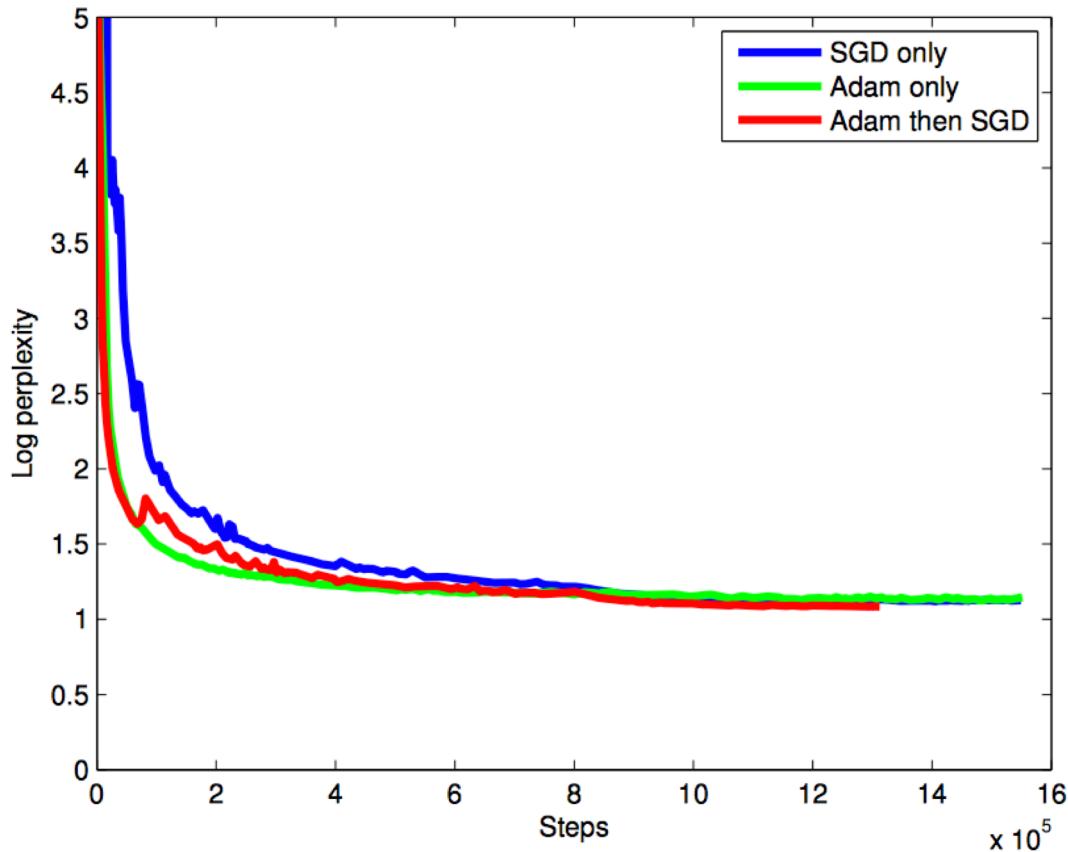
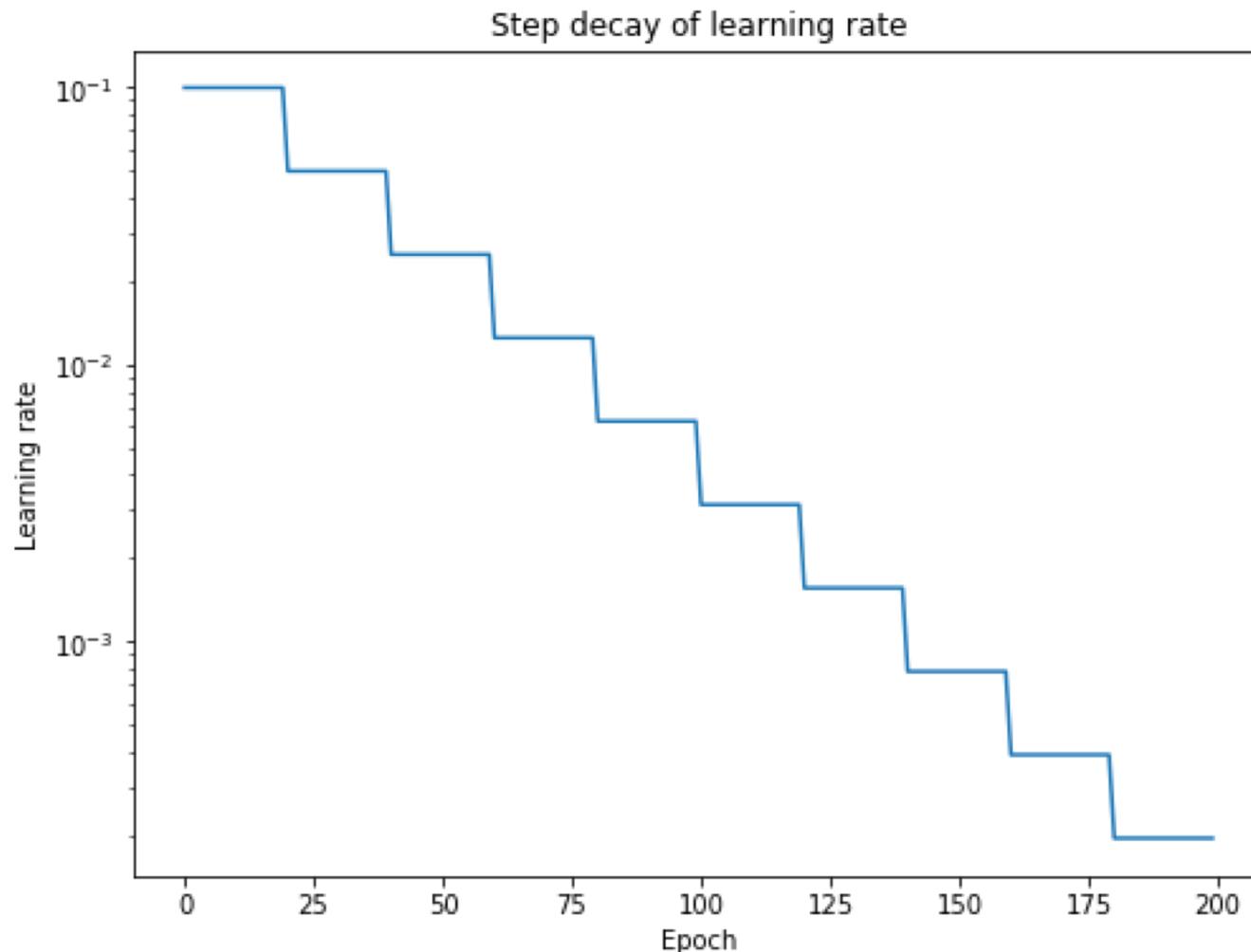
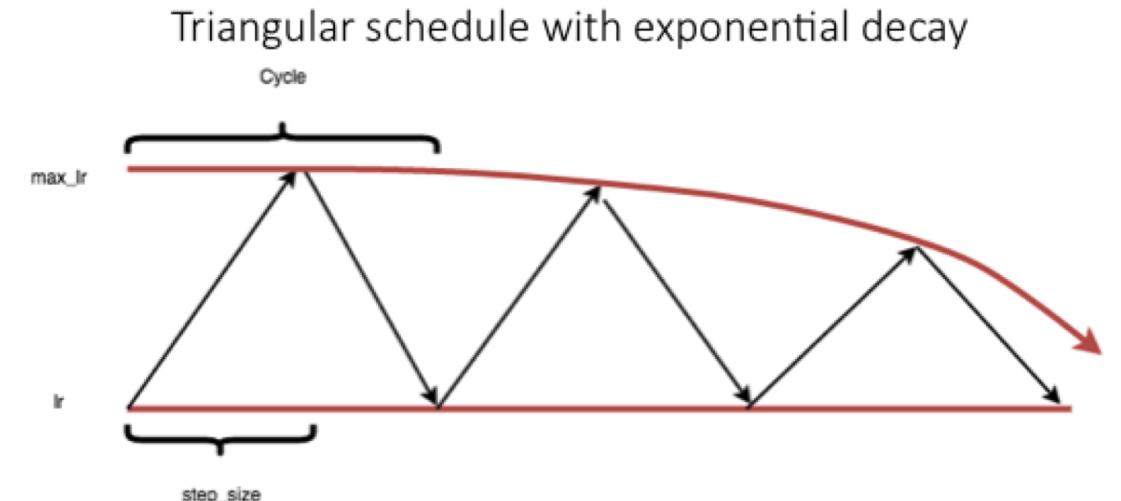
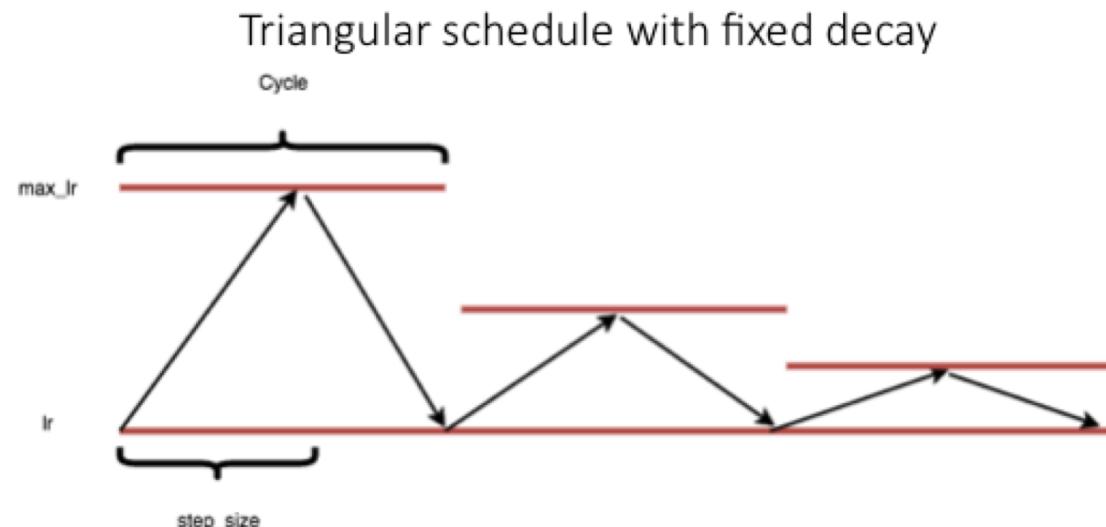
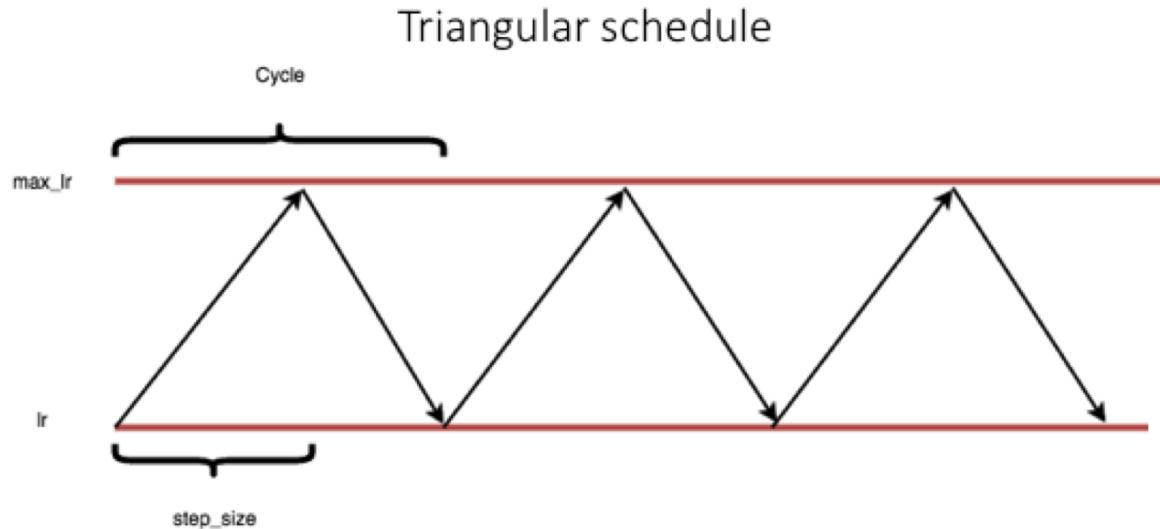


Figure 5: Log perplexity vs. steps for Adam, SGD and Adam-then-SGD on WMT En→Fr during maximum likelihood training. Adam converges much faster than SGD at the beginning. Towards the end, however, Adam-then-SGD is gradually better. Notice the bump in the red curve (Adam-then-SGD) at around 60k steps where we switch from Adam to SGD. We suspect that this bump occurs due to different optimization trajectories of Adam vs. SGD. When we switch from Adam to SGD, the model first suffers a little, but is able to quickly recover afterwards.

Learning Rate Annealing



Cyclic Learning Rate



(Smith, 2015)

Cyclic Learning Rate

We can write the general schedule as

$$\eta_t = \eta_{\min} + (\eta_{\max} - \eta_{\min}) (\max (0, 1 - x))$$

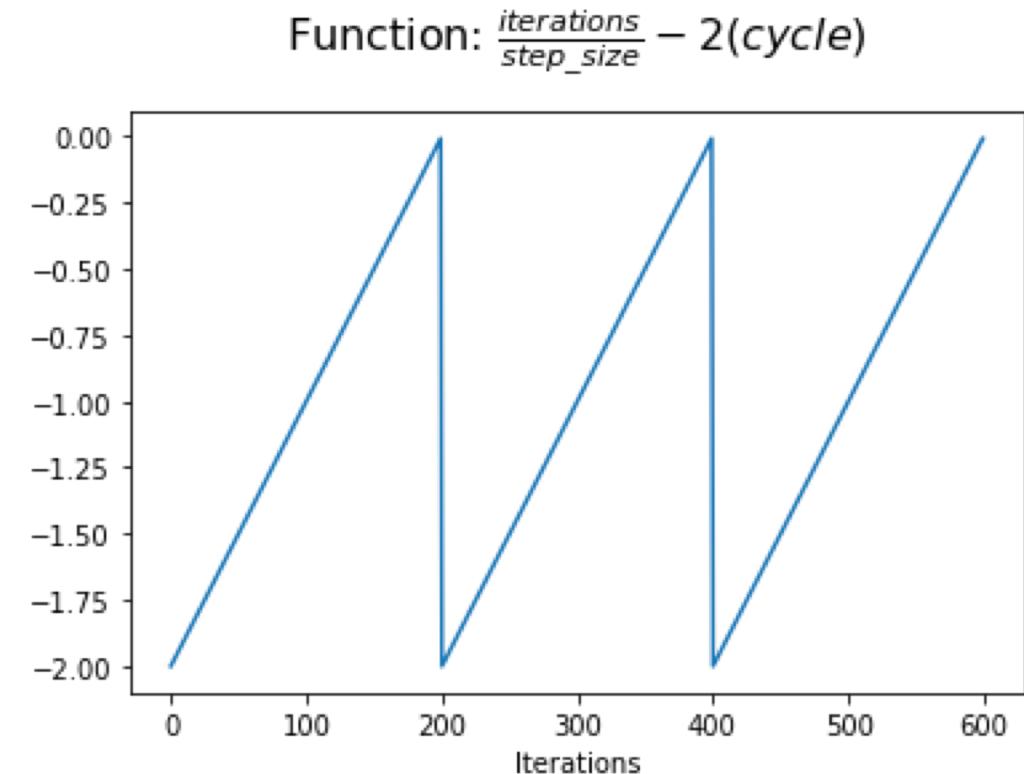
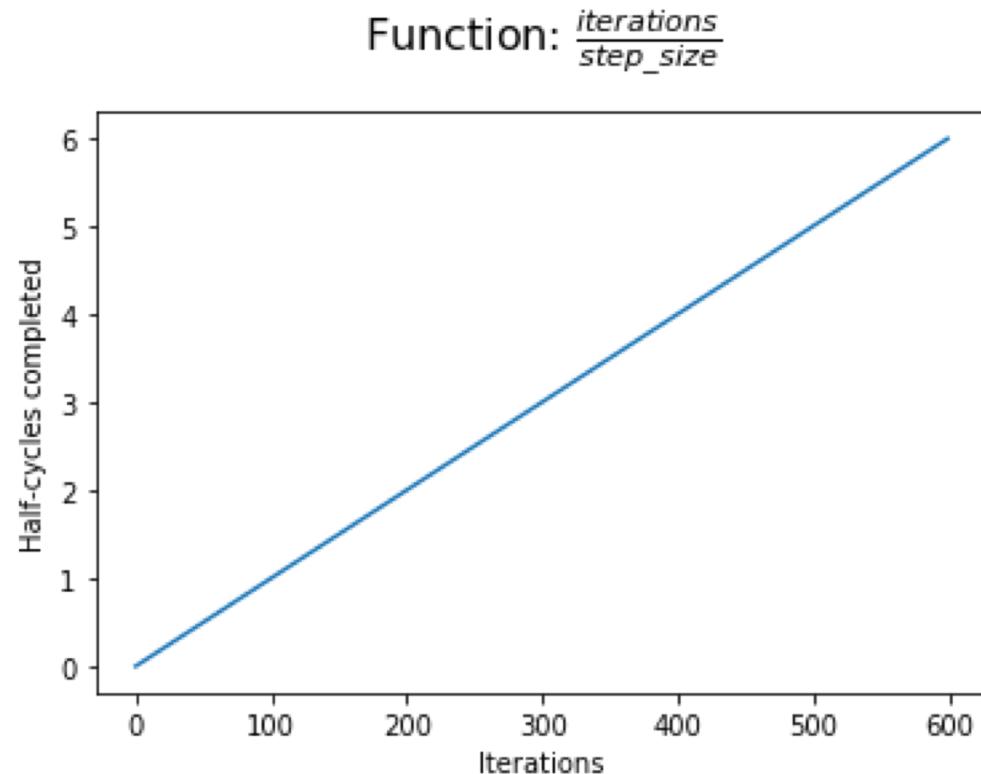
where x is defined as

$$x = \left| \frac{\text{iterations}}{\text{stepsize}} - 2(\text{cycle}) + 1 \right|$$

and cycle can be calculated as

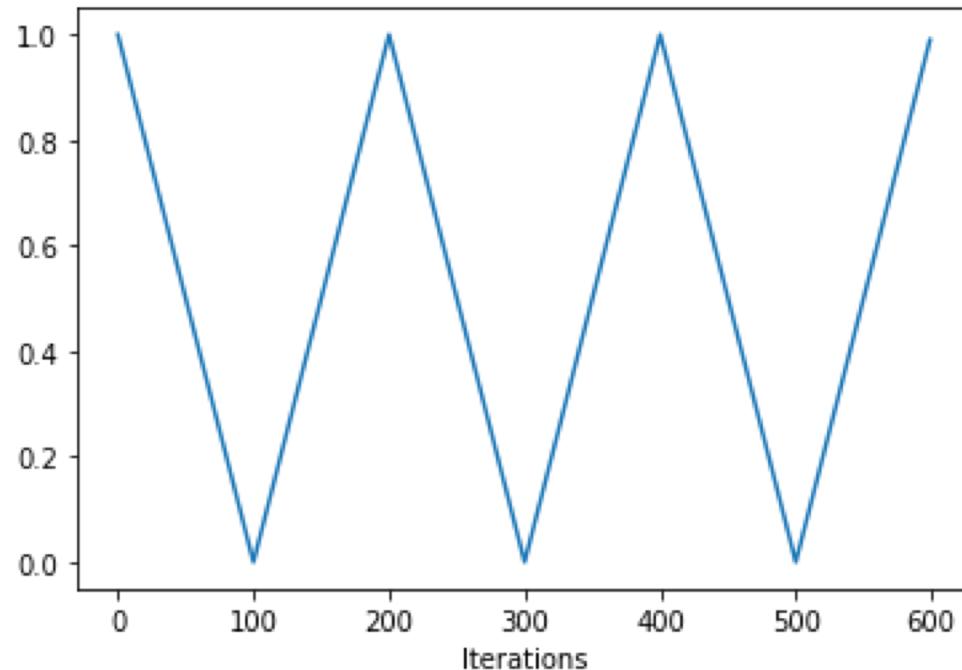
$$\text{cycle} = \text{floor} \left(\frac{1 + \text{iterations}}{2(\text{stepsize})} \right)$$

Cyclic Learning Rate

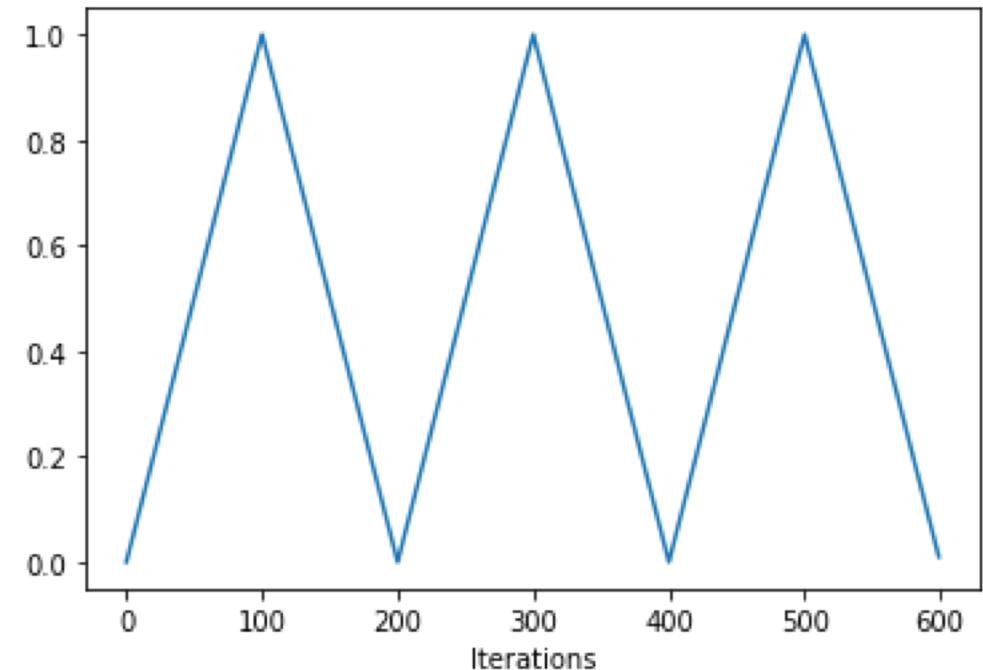


Cyclic Learning Rate

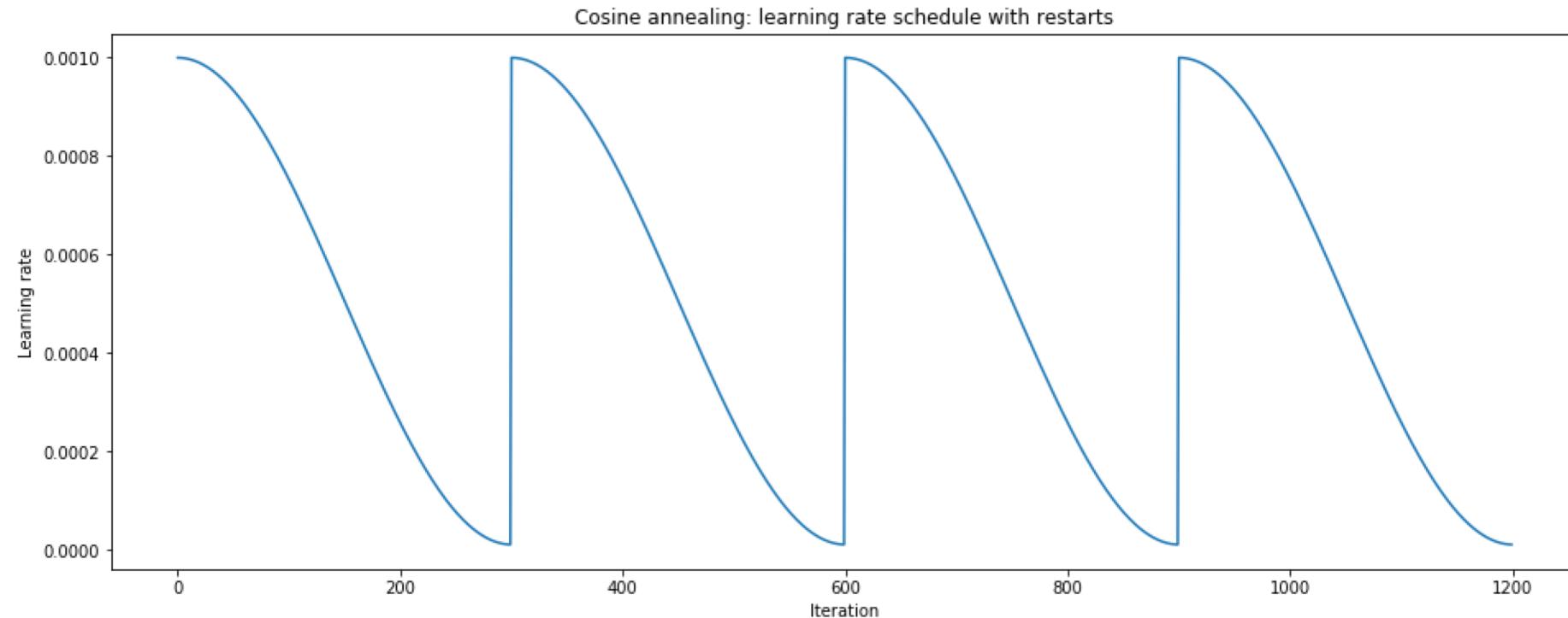
Function: $\left\lfloor \frac{\text{iterations}}{\text{step_size}} - 2(\text{cycle}) + 1 \right\rfloor$



Function: $1 - \left\lfloor \frac{\text{iterations}}{\text{step_size}} - 2(\text{cycle}) + 1 \right\rfloor$



Cosine SGD with Warm Restarts



$$\eta_t = \eta_{\min}^i + \frac{1}{2} (\eta_{\max}^i - \eta_{\min}^i) \left(1 + \cos\left(\frac{T_{current}}{T_i} \pi\right) \right)$$

Annealing,
i.e. taking a
partial step

```

2 import torch
3 import torch.nn as nn
4 import torch.nn.functional as F
5 import torch.optim as optim
6
7 from CLR_preview import CyclicLR
8 from adamW import AdamW
9
10 # Step 1: Initialization.
11 model = nn.Sequential(
12     nn.Linear(input_dim, hidden_dim),
13     nn.Sigmoid(),
14     nn.Linear(hidden_dim, output_dim),
15     nn.Sigmoid()
16 )
17
18 optimizer = AdamW(model.parameters(), lr=0.001, betas=(0.9, 0.99), weight_decay = 0.1)
19 # For classification.
20 ##clr_stepsize = (num_classes*50//int(batchsz))*4
21 clr_stepsize = 3e-4
22 clr_wrapper = CyclicLR(optimizer, step_size=clr_stepsize)
23
24 # Initialize the loss function.
25 criterion = nn.MSELoss()
26
27 losses = [] # Keeps track of the losses.
28 # Step 2-4 of training routine.
29 for _e in tqdm(range(num_epochs)):
30     # Reset the gradient after every epoch.
31     optimizer.zero_grad()
32     # Step 2: Foward Propagation
33     predictions = model(X)
34     # Step 3: Back Propagation
35     # Calculate the cost between the predictions and the truth.
36     loss = criterion(predictions, Y)
37     # Remember to back propagate the loss you've computed above.
38     loss.backward()
39     # Step 4: Optimizer take a step and update the weights.
40     optimizer.step()
41     losses.append(loss.data.item())
42

```

Super convergence, see
<https://www.fast.ai/2018/07/02/ada-m-weight-decay/>

CLR_preview from
https://github.com/ahirner/pytorch-retraining/blob/master/CLR_preview.py

AdamW from
<https://github.com/egg-west/AdamW-pytorch>

Environment Setup

Open Anaconda Navigator.

Go to the PyTorch installation page, copy the command as per configuration:

<https://pytorch.org/get-started/locally/>

Fire up the terminal in Anaconda Navigator.

Start a Jupyter Notebook.

Download <http://bit.ly/ANLP-Session4>

Import the .ipynb to the Jupyter Notebook

Fin